



# Assessment of individual tumor buds using keratin immunohistochemistry: moderate interobserver agreement suggests a role for machine learning

J. M. Bokhorst<sup>1</sup> · A. Blank<sup>2</sup> · A. Lugli<sup>2</sup> · I. Zlobec<sup>2</sup> · H. Dawson<sup>2</sup> · M. Vieth<sup>3</sup> · L. L. Rijstbergen<sup>1</sup> · S. Brockmoeller<sup>4</sup> · M. Urbanowicz<sup>5</sup> · J. F. Flejou<sup>6</sup> · R. Kirsch<sup>7</sup> · F. Ciampi<sup>1</sup> · J. A. W. M. van der Laak<sup>1,8</sup> · I. D. Nagtegaal<sup>1</sup>

Received: 20 September 2019 / Revised: 7 November 2019 / Accepted: 23 November 2019

© The Author(s) 2019, corrected publication January 2020

## Abstract

Tumor budding is a promising and cost-effective biomarker with strong prognostic value in colorectal cancer. However, challenges related to interobserver variability persist. Such variability may be reduced by immunohistochemistry and computer-aided tumor bud selection. Development of computer algorithms for this purpose requires unequivocal examples of individual tumor buds. As such, we undertook a large-scale, international, and digital observer study on individual tumor bud assessment. From a pool of 46 colorectal cancer cases with tumor budding, 3000 tumor bud candidates were selected, largely based on digital image analysis algorithms. For each candidate bud, an image patch (size 256 × 256 μm) was extracted from a pan cytokeratin-stained whole-slide image. Members of an International Tumor Budding Consortium ( $n = 7$ ) were asked to categorize each candidate as either (1) tumor bud, (2) poorly differentiated cluster, or (3) neither, based on current definitions. Agreement was assessed with Cohen's and Fleiss Kappa statistics. Fleiss Kappa showed moderate overall agreement between observers (0.42 and 0.51), while Cohen's Kappas ranged from 0.25 to 0.63. Complete agreement by all seven observers was present for only 34% of the 3000 tumor bud candidates, while 59% of the candidates were agreed on by at least five of the seven observers. Despite reports of moderate-to-substantial agreement with respect to tumor budding grade, agreement with respect to individual pan cytokeratin-stained tumor buds is moderate at most. A machine learning approach may prove especially useful for a more robust assessment of individual tumor buds.

## Introduction

Tumor budding, defined as the presence of single cells or small clusters of up to four tumor cells, may be seen at the invasive front of colorectal cancer and other tumor types [1]. Tumor buds (TB) are detached (epithelial) tumor cells that, in close interaction with their microenvironment, transform at least partially into a mesenchymal stem-like phenotype. In the process, TB lose some of their epithelial characteristics and acquire features that are correlated with increased cell motility [2]. TB are morphologically and biologically related to poorly differentiated clusters (PDC's), which represent larger tumor cell clusters (five or more cells without gland formation), and which are also located at the tumor invasive front [3].

The clinical significance of the extent of TB as an independent risk factor for adverse outcomes in colorectal cancer (CRC) is now well established [4–8]. However, until recently the reporting of tumor budding in routine practice

The original version of this article was revised due to a retrospective Open Access order.

✉ J. M. Bokhorst  
john-melle.bokhorst@radboudumc.nl

<sup>1</sup> Radboud University Medical Center, Nijmegen, Netherlands

<sup>2</sup> University of Bern, Bern, Switzerland

<sup>3</sup> University of Bayreuth, Bayreuth, Germany

<sup>4</sup> University of Leeds, Leeds, UK

<sup>5</sup> EORTC Translational Research Unit, Brussels, Belgium

<sup>6</sup> Saint-Antoine Hospital, Paris, France

<sup>7</sup> University of Toronto, Toronto, Canada

<sup>8</sup> Center for Medical Image Science and Visualization, Linköping University, Linköping, Sweden

was held back by a lack of internationally accepted criteria and methodology for its assessment. The recent publication of evidence-based recommendations for TB assessment, formulated during the International Tumor Budding Consensus Conference (ITBCC, 2016), removed an important barrier to the wider reporting of TB [1].

As with all adverse histologic features (and especially those requiring quantitation), interobserver agreement of tumor budding is suboptimal. The degree of interobserver variability may be influenced by the type of stain used (H&E vs. immunohistochemistry), the number of output categories, cut-off values, and experience/expertise of the observers [9–11]. Although reported concordance between observers has generally been within acceptable limits, ongoing concerns regarding interobserver variability have continued to hamper clinical adoption. Interobserver variability in TB assessment may be greatly reduced by computer-aided detection systems. A class of algorithms based on deep learning and capable of recognizing complex tissue structures and patterns is emerging as a promising tool in medical imaging and digital pathology, matching, and in some cases surpassing human performance [12–14]. A requirement for successful development of such supervised deep learning algorithms is the availability of a sufficient amount of correctly labeled input data. As such, its application to TB assessment requires a large number of images containing TB which have been confirmed by (a panel of) experienced pathologists.

We have previously assessed tumor budding [15] in digital images from preselected hotspot locations, both on H&E and in cytokeratin 8–18 (CK) stained slides with the purpose of identifying TB for training computer algorithms. Moderate Kappa scores were found between two pathologists in scoring tumor budding grade, and notably few individual buds were annotated by both pathologists (H&E 38% and CK 54% of the total annotated TB), indicating considerable interobserver variation in identifying individual TB.

In the present study, we aim to establish the level of agreement among expert pathologists when assessing a series of 3000 individual tumor bud candidates in digital images of CK stained slides. Our study differs from other observer studies, which have largely focused on overall tumor budding grade rather than individual TB. Digitized images of 3000 TB-candidates were submitted to an international panel of experts, with the request to determine per candidate status as follows: TB, PDC, or neither. We have chosen to derive these TB/PDC candidates mainly from pan-CK stained tissue. For comparison purposes 150 of the 3000 TB-candidates were also offered for evaluation in ITBCC preferred H&E stained version separately.

## Material and methods

### Materials

From three centers, a total of 46 CRC patients were selected, in which routine diagnostic assessment had revealed the presence of TB. Per patient, one paraffin-embedded tissue block was included. Two slides were stained with AE1–AE3 immunohistochemistry at the Dublin University Hospital, five slides were stained with AE1–AE3 immunohistochemistry at Bern University Labs. The remaining 39 slides were stained with H&E, digitally captured (producing whole slide images; WSI), subsequently destained, restained with CK8-18 immunohistochemistry, and scanned again at the Radboud UMC in Nijmegen (procedure described by van den Brand et al. [16]). All slides were scanned with a Panoramic P250 Flash II scanner (3D-Histech, Hungary) using a 40× objective lens (yielding specimen level pixel-size of  $0.24 \times 0.24 \mu\text{m}$ ). As only fully anonymized archival tissue was used in this study, the need for ethical approval was waived by the institutional review board of Radboud UMC.

The set of 3000 TB-candidates was comprised of 800 TB-candidates, taken from a previous study [15]. The remaining 2200 TB-candidates were selected as follows. One pathologist marked hotspot regions (measuring  $0.785 \text{ mm}^2$ ) in the immunohistochemistry stained slides. Next, an image analysis algorithm was used to automatically identify individual bud candidates. The algorithm performed a color deconvolution [17] to isolate DAB positive image pixels, which were grouped to form binary objects. Only objects with a surface area between 25 and  $5000 \mu\text{m}^2$  were retained, to remove small artifacts and larger clusters of tumor cells. Candidates were visually checked by an expert and obvious artifacts were omitted. From the remaining objects, 1900 TB-candidates with surface area  $<1000 \mu\text{m}^2$  were taken (mostly representing TB) and another 300 candidates with surface area  $>1000 \mu\text{m}^2$  (containing mainly PDC's) were randomly chosen, to ensure sufficient presence of TB. The set of 3000 candidates was randomly split into two groups of 1500 cases.

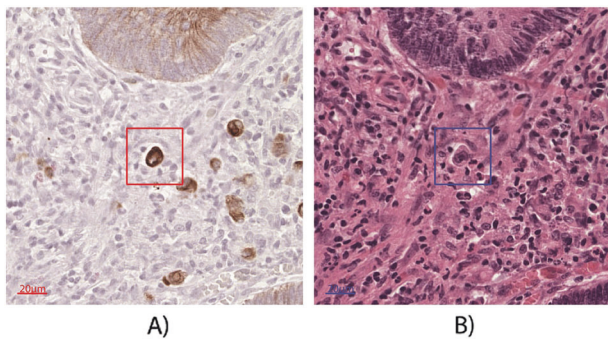
Restaining of the slides allowed us to select identical objects in both H&E and CK immunohistochemistry stains. A part of the 3000 candidates was traced based on their immunohistochemistry image coordinates into the equivalent images in H&E stain. A total of 150 TB-candidates were sampled in H&E, selecting these randomly from different groups (i.e., TB, PDC, neither) in proportion to the number of objects per group.

### Method

Eleven pathologists from seven different institutions across seven different countries, all with expertise in the field of

tumor budding, were invited to participate in this study. Pathologists were asked to evaluate two study sets each including 1500 TB/PDC candidates in immunohistochemistry and, 2 months later, a set of 150 TB/PDC candidates on H&E according to ITBCC guidelines and best practice, including the rule to regard a visible presence of the nucleus as a requirement. Observers were asked to designate the candidates into one of three categories, namely (1) TB, (2) PDC, or (3) neither (no TB, no PDC). The “gold standard” status of an object was defined using a voting rule of 70% or more: the label is assigned if at least 5 out of 7 and 8 out of 11 pathologists agree in immunohistochemistry and H&E, respectively. To enable simultaneous execution of these tasks by different pathologists from the various work locations, a web-based platform was used. Through this platform, patches of  $0.25 \times 0.25$  mm were shown, based on the center of mass (CoM) of each candidate. To clearly mark the object of interest, a square with a fixed size of  $0.03 \text{ mm}^2$  was added, based again on its CoM. An example is shown in Fig. 1.

Fleiss Kappa and Cohen’s Kappa were used to calculate agreement between multiple observers and different combinations of two observers (one vs. one) respectively. For the Kappa scores we used the terms formulated by Landis and Koch [18], i.e., <0.2 poor, 0.21–0.40 fair, 0.40–0.60 moderate, 0.61–0.80 good, and >0.80 very good.



**Fig. 1** Example of TB-candidate in IHC and re-stained H&E. Example of TB-candidate in **a** immunohistochemistry, and **b** restained H&E.

## Results

Immunohistochemistry results were obtained from nine pathologists, five of whom scored both sets of 1500 TB/PDC candidates. The remaining four pathologists scored either the first or the second 1500 candidates. As a result, 7 scores were obtained per TB/PDC object.

Of the 3000 TB/PDC candidates, 612, 15, and 386 cases were unanimously classified as TB positive, PDC positive, or neither, respectively. Based on the previously agreed 70%-majority vote, we came to a further definite classification for 1765 candidates (1010 TB, 52 PDC, and 703 neither). For the remaining 1235 candidates, no consensus could be achieved. An overview of all scores can be found in Table 1; examples of every category can be found in Fig. 2.

Of the  $2 \times 1500$  immunohistochemistry candidates, on average  $803 (\pm 194)$ ,  $67 (\pm 28)$ , and  $630 (\pm 201)$  were classified as TB, PDC, or neither in Group 1. In Group 2,  $758 (\pm 118)$ ,  $111 (\pm 31)$ , and  $631 (\pm 143)$  candidates were labeled as TB, PDC, or neither (Table 2). Individual scores are shown in Fig. 3.

Fleiss Kappa was calculated for the first 1500 TB-candidates (classification result Group 1) and the second 1500 TB-candidates (classification result Group 2) separately. With values of 0.42 for the first group and 0.51 for the second group, assessment in both groups showed moderate agreement.

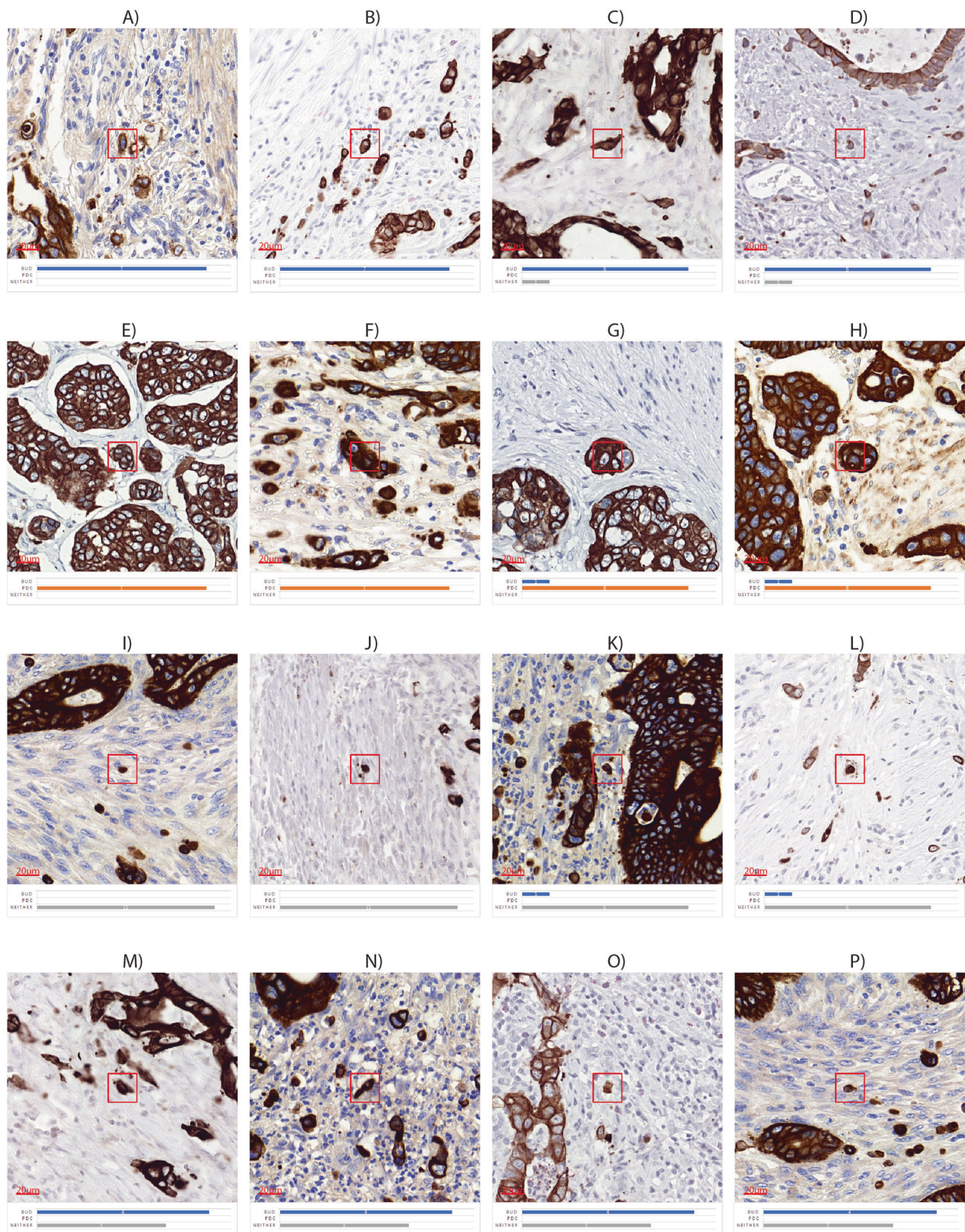
One-versus-one observer agreement numbers (Cohen’s Kappa) are shown per group in Fig. 4. In the first group these scores vary between 0.24 and 0.60, in the second group they range from 0.36 to 0.65. In the second observer group, the 2nd, 3rd, 4th, and 5th observer achieved a higher agreement score compared with each other. This is also reflected in the lower standard deviation numbers of the second group compared with the first one.

A series of 150 TB/PDC candidates on H&E stained sections was assessed by 11 observers. A Fleiss Kappa of 0.49 was achieved on this subset. The Cohen Kappa numbers ranged from 0.07 to 0.782 (Fig. 5). Of the 150 H&E candidates on average  $54 (\pm 18)$ ,  $13 (\pm 2)$ , and  $84 (\pm 19)$ , were classified as TB, PDC, or neither. Individual scores can be found in Fig. 6.

**Table 1** Number of objects, classified as TB, PDC, or neither by all, majority, or minority (no agreement) per candidate group (3000 IHC, 150 H&E, and 150 IHC).

	Uniform			70% majority			No agreement (all classes)
	Bud	PDC	Neither	Bud	PDC	Neither	
IHC (3000 objects)	612 (20.40%)	15 (0.50%)	386 (12.87%)	398 (13.27%)	37 (1.23%)	317 (10.57%)	1235 (41.17%)
H&E (150 objects)	8 (5.33%)	1 (0.67%)	13 (8.67%)	31 (20.67%)	10 (6.67%)	68 (45.33%)	19 (12.67%)
IHC reference group (150 objects)	26 (17.33%)	5 (3.33%)	25 (16.00%)	25 (16.67%)	7 (4.67%)	24 (16.00%)	38 (25.33%)





**Fig. 2** Examples of TB-candidates with manual scores. Examples of TB-candidates **a, b** uniform selected as bud, **c, d** majority vote bud, **e, f** uniform PDC, **g, h** majority vote PDC, **i, j** uniform neither,

**k, l** majority vote neither. **m–p** no agreement was reached. Legend of colors: blue—bud, orange—PDC, gray—neither.

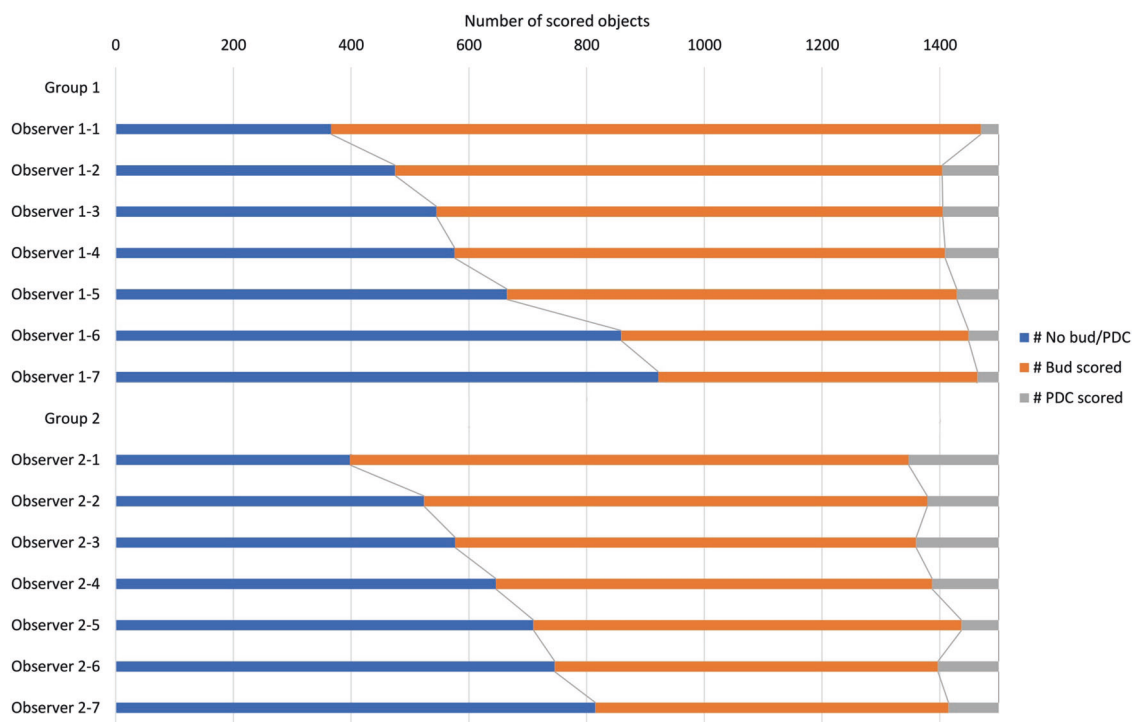
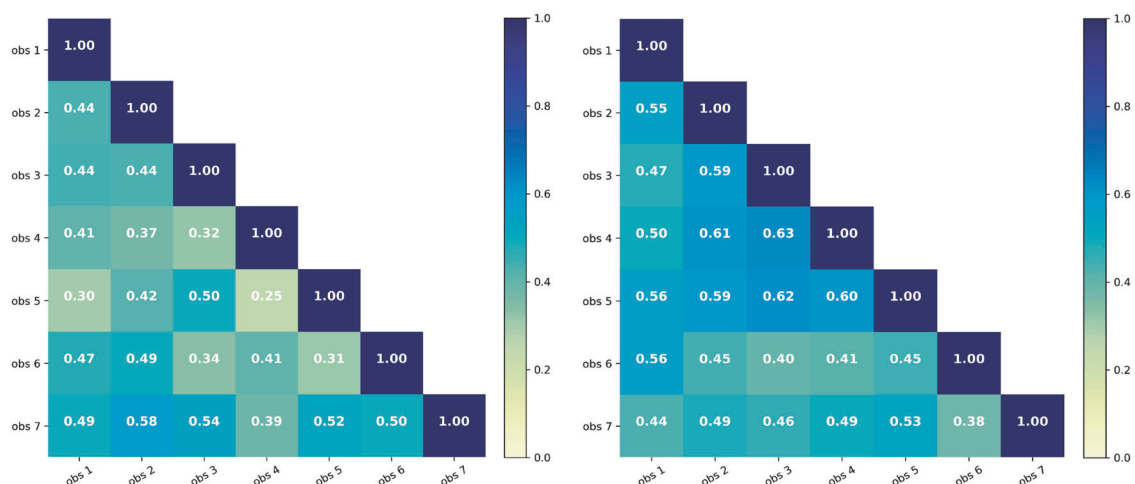
**Table 2** Averaged number of objects (with std. dev.) per observer group.

Per observer group × dataset: Mean number (std. dev.) of TB, PDC, or neither

	Bud	PDC	Neither
IHC—Group 1 (1500 objects)	803 (193)	67 (28)	630 (201)
IHC—Group 2 (1500 objects)	758 (118)	111 (30)	631 (142)
H&E (150 objects)	54 (18)	13 (2)	84 (19)
IHC reference group (150 objects)	70 (14)	12 (2)	68 (14)

Table 1 shows that 8, 1, and 13 cases were unanimously classified as TB positive, PDC positive, and neither respectively. Based on the majority vote we counted 39 TB, 11 PDC, and 81 neither. On 19 cases no agreement could be reached.

The number of “no agreement” objects was therefore decreased in H&E compared with the immunohistochemistry reference group. Numbers of PDC remained virtually unchanged, but the number of objects, classified as TB by at least 5 of the 7 observers (Uniform + Majority), decreased on balance with higher numbers of neither class (25 + 24 versus

**Fig. 3** TB, PDC and neither scores in the 2 × 1500 immunohistochemistry dataset per observer.**Fig. 4** One-versus-one Cohen Kappa scores per observer of left group 1 and right group 2.



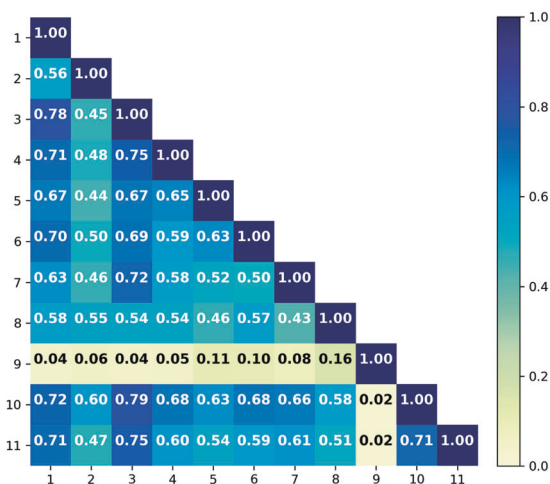
13 + 68 objects). Table 2 also shows that the average number of PDCs in the H&E group almost corresponds to the average number of PDCs in the immunohistochemistry reference group. Example cases can be found in Fig. 7.

## Discussion

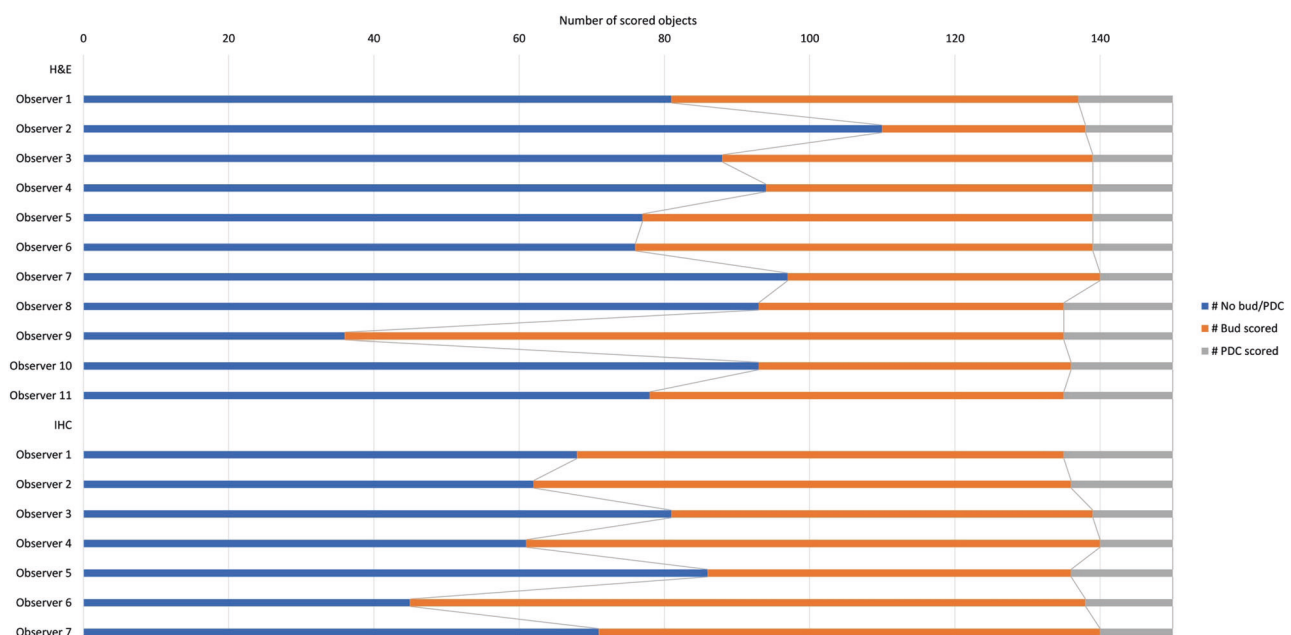
In contrast to previous studies, we assessed observer variability in scoring individual TB, rather than the overall budding. This approach prevented pathologists from unsubconsciously taking the tumor grade and differentiation grade into account when assessing individual TB. Accurate

identification of individual TB is critical for subsequent development of artificial intelligence models which can automate this task, thereby potentially reducing observer bias. We found that the agreement between pathologists for scoring buds in immunohistochemistry was only moderate, which was slightly improved if buds were identified in H&E. Only one in three of the 3000 immunohistochemically stained candidates was uniformly classified as either TB, PDC, or neither by unanimity or majority vote. No agreement was reached on average for four out of ten candidates. Since interobserver variability in the assessment of individual TB can result in variability in tumor budding grades, it is of interest to investigate the cause(s) of such variability.

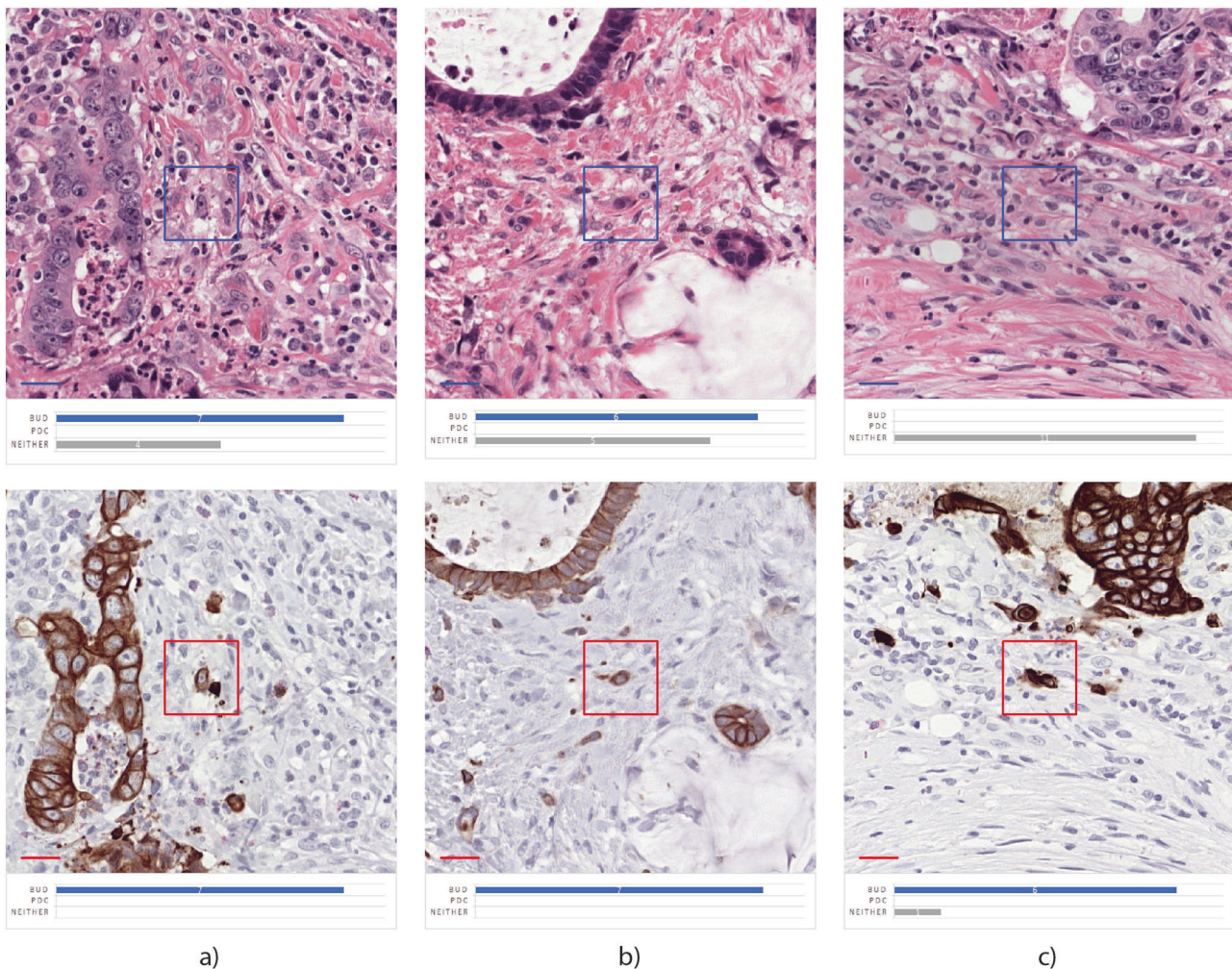
In the interpretation and assessment of structures in immunohistochemistry, TB/PDCs can be mistaken for other cytokeratin-positive objects. These objects can be partially described as remnants, parts of epithelium that is at some stage of degradation. Various studies provide examples of this. Mitrovic et al. [19] cites fragmentation of neoplastic glandular structures as a possible cause of misclassification. In line with this, Kai et al. [20] call the presence of “pseudo buds” an unexpected pitfall in the assessment of TB in immunohistochemistry, especially in cases with severe inflammation. They consider pseudo buds to be comprised of fragmented epithelium, destroyed by inflammatory cells. Koelzer et al. [9] note that the assessment of TB in immunohistochemistry can be complicated by the presence of cytokeratin-positive microvesicles and membrane fragments. Although not



**Fig. 5** Cohen Kappa scores per observer for the 150 H&E cases.



**Fig. 6** TB, PDC, and neither scores in the 150 H&E dataset per observer with the addition of scores of the identical objects in immunohistochemistry staining.



**Fig. 7** Score shifted examples of restrained TB-candidates. **a, b** Shift in the agreement from uniform bud vote in immunohistochemistry to no agreement in H&E, **c** shift from majority vote on tumor bud in immunohistochemistry to uniform neither assessment in H&E.

mentioned, degenerating or apoptotic cells are also part of this CK positive group.

In contrast, Shinto et al. [21] researched the morphology, configuration, and role of small nonnucleated cytoplasmic fragments (minimum diameter 2  $\mu\text{m}$ ) often detected around TB. The authors observe that some cytoplasmic fragments are clearly connected with budding foci, suggesting that they represent dendritic cell processes, rather than isolated cell fragments. They, therefore, have renamed them cytoplasmic pseudo fragments.

As we used a minimum area of 25  $\mu\text{m}^2$  for the computer-aided preselection of TB-candidates, cytoplasmic pseudo fragments were removed a priori from the candidate group. As a consequence, interobserver variability can only be associated with the first-mentioned objects in this study. Since no definite cause could be found in the candidates for whom no agreement was reached, further research will have to show whether/to what extent each of them actually contributes to interobserver disagreement.

Clear presence of a cell nucleus is an important criterion for TB/PDC as it allows observers to distinguish TB/PDC from aforementioned CK positives.

Although CK staining helps in the first phase of detection of the (malignant) epithelium, e.g., the hotspot selection, slight overstaining can be counterproductive because informative variations in intensity that help the interpretation of morphological structures can be lost for the human observer. It is therefore possible that part of the disagreement in immunohistochemistry can be traced to poor visibility of the nuclei in this study, where observers did not have much opportunity to fall back on relevant information from the context, i.e., nearby gland fully intact or not, position of the candidate relative to the gland, but also frame of reference for estimating area/size of the candidate.

As is well known, the human eye is quite poor at judging minor variations in intensity, but computers are especially good at assessing finer shades of intensity. A number of researchers have now entered the field of automatic

detection of TB in IHC stained slides. In a recent article, Bergler et al. [22] present a hybrid method. In a first phase detection step, TB-candidates are generated using classical image processing methods, whereby the segmentation of pan-CK positive objects primarily is performed on color information or signal intensity thresholds and information on size. A second operation then is added to this phase, which is aimed at reducing the number of false positives and carried out with the help of deep learning-based algorithms. Weis et al. [23] follow a similar protocol for immunohistochemistry stained TMA cores. Both authors show promising results but need further validation before implementation in daily clinical routine.

When H&E sections were presented, a group of 11 observers achieved the same moderate overall performance on the assessment of 150 TB/PDC candidates. The inter-observer agreement on H&E is similar to the level of agreement on immunohistochemistry. We did see a shift in the classifications, however, there was a decrease in candidates classified as TB and an increase in candidates classified as neither. The variability in classification of TB's on H&E—other than in immunohistochemistry—may be related to challenges in differentiating TB from inflammatory cells (lymphocytes, histiocytes, and macrophages) and stromal cells, as activated fibroblasts [19]. Based on this observation, we hypothesize that interobserver variability may be less attributable to TB composed of clusters of 2–4 cells than to individual cell TB. As the number of cells of the candidates was not registered, we could not test this hypothesis in this study. It is worth considering this in future investigations.

Relevant morphological cell characteristics are better preserved in H&E, but TB detection is then more time-consuming. Furthermore, there is a greater chance of TB's being overlooked compared to IHC, where the main issue is “overcall” due to the presence of TB mimics. Nevertheless, based on the results of this study, there is no reason to deviate from the ITBCC preference for H&E staining at this time. It should be noted, however, that only a relatively small number of TB-candidates were classified in H&E sections.

In conclusion, we have shown that the assessment of individual TB-candidates in immunohistochemistry is difficult, even for pathologists with expertise in the field of TB. Although immunohistochemical staining helps with the detection of TB/PDC candidates, there is a risk of misclassification in connection with the common presence of cytokeratin-positive “competitors”. The assessment of these objects must be done correctly, on the basis of nuclear absence. Doubtful (in-)visibility of the cell nucleus in immunohistochemistry is in this context a complicating factor that can be remedied with the use of computers.

In this study we have already applied a form of computer-aided selection by automated preselection based on color and size. With the results obtained, further steps can now be taken to achieve automation of TB assessment.

**Acknowledgements** This project has received funding from the Dutch Cancer Society, project number 10602/2016-2. The authors would like to thank Dr David Gibbons for his participation in scoring the H&E candidates. The authors would also like to thanks Femke Doubrava-Simmer for general assistance and many discussions.

## Compliance with ethical standards

**Conflict of interest** The authors declare that they have no conflict of interest.

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

1. Lugli A, Kirsch R, Ajioka Y, Bosman F, Cathomas G, Dawson H, et al. Recommendations for reporting tumor budding in colorectal cancer based on the International Tumor Budding Consensus Conference (ITBCC) 2016. *Mod Pathol*. 2017;30:1299–311.
2. Gabbert H, Wagner R, Moll R, Gerharz C-D. Tumor dedifferentiation: an important step in tumor invasion. *Clin Exp Metastasis*. 1985;3:257–79.
3. Hong M, Kim JW, Shin MK, Kim BC. Poorly differentiated clusters in colorectal adenocarcinomas share biological similarities with micropapillary patterns as well as tumor buds. *J Korean Med Sci*. 2017;32:1595–602.
4. Kawachi H, Eishi Y, Ueno H, Nemoto T, Fujimori T, Iwashita A, et al. A three-tier classification system based on the depth of submucosal invasion and budding/sprouting can improve the treatment strategy for T1 colorectal cancer: a retrospective multicenter study. *Mod Pathol*. 2015;28:872.
5. Ueno H, Murphy J, Jass J, Mochizuki H, Talbot I. Tumour budding as an index to estimate the potential of aggressiveness in rectal cancer. *Histopathology*. 2002;40:127–32.
6. Park K-J, Choi H-J, Roh M-S, Kwon H-C, Kim C. Intensity of tumor budding and its prognostic implications in invasive colon carcinoma. *Dis Colon Rectum*. 2005;48:1597–602.
7. Losi L, Ponti G, Di Gregorio C, Marino M, Rossi G, Pedroni M, et al. Prognostic significance of histological features and biological parameters in stage I (pT1 and pT2) colorectal adenocarcinoma. *Pathol Res Pract*. 2006;202:663–70.



8. Yamauchi H, Togashi K, Kawamura YJ, Horie H, Sasaki J, Tsujinaka S, et al. Pathological predictors for lymph node metastasis in T1 colorectal cancer. *Surg Today*. 2008;38:905–10.
9. Koelzer VH, Zlobec I, Berger MD, Cathomas G, Dawson H, Dirschmid K, et al. Tumor budding in colorectal cancer revisited: results of a multicenter interobserver study. *Virchows Arch*. 2015;466:485–93.
10. Koelzer VH, Zlobec I, Lugli A. Tumor budding in colorectal cancer—ready for diagnostic practice? *Hum Pathol*. 2016;47:4–19.
11. Martin B, Schafer E, Jakubowicz E, Mayr P, Ihringer R, Anthuber M, et al. Interobserver variability in the H&E-based assessment of tumor budding in pT3/4 colon cancer: does it affect the prognostic relevance? *Virchows Arch*. 2018;473:189–97.
12. Bejnordi BE, Veta M, Van Diest PJ, Van Ginneken B, Karssemeijer N, Litjens G, et al. Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer. *JAMA*. 2017;318:2199–210.
13. Nagpal K, Foote D, Liu Y, Chen P-HC, Wulczyn E, Tan F, et al. Development and validation of a deep learning algorithm for improving Gleason scoring of prostate cancer. *npj Digit Med*. 2019;2:48.
14. Campanella G, Hanna MG, Geneslaw L, Miraflor A, Silva VWK, Busam KJ, et al. Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. *Nat Med*. 2019;25:1301–9.
15. Bokhorst J-M, Rijstenberg L, Goudkade D, Nagtegaal I, van der Laak J, Ciompi F. Automatic detection of tumor budding in colorectal carcinoma with deep learning. *Computational Pathology and Ophthalmic Medical Image Analysis: First International Workshop, COMPAY 2018, and 5th International Workshop, OMIA 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16–20, 2018*. In: *Computational pathology and ophthalmic medical image analysis*. Springer; 2018. p. 130–8.
16. van den Brand M, Hoevenaars BM, Sigmans JH, Meijer JW, van Cleef PH, Groenen PJ, et al. Sequential immunohistochemistry: a promising new tool for the pathology laboratory. *Histopathology*. 2014;65:651–7.
17. Ruifrok AC, Johnston DA. Quantification of histochemical staining by color deconvolution. *Anal Quant Cytol Histol*. 2001;23:291–9.
18. Landis JR, Koch GG. The measurement of observer agreement for categorical data. *Biometrics*. 1977;33:159–74.
19. Mitrovic B, Schaeffer DF, Riddell RH, Kirsch R. Tumor budding in colorectal carcinoma: time to take notice. *Mod Pathol*. 2012;25:1315–25.
20. Kai K, Aishima S, Aoki S, Takase Y, Uchihashi K, Masuda M, et al. Cytokeratin immunohistochemistry improves interobserver variability between unskilled pathologists in the evaluation of tumor budding in T1 colorectal cancer. *Pathol Int*. 2016;66:75–82.
21. Shinto E, Mochizuki H, Ueno H, Matsubara O, Jass J. A novel classification of tumour budding in colorectal cancer based on the presence of cytoplasmic pseudo-fragments around budding foci. *Histopathology*. 2005;47:25–31.
22. Bergler M, Benz M, Rauber D, Hartmann D, Kötter M, Eckstein, et al. Automatic Detection of Tumor Buds in Pan-Cytokeratin Stained Colorectal Cancer Sections by a Hybrid Image Analysis Approach. In *European Congress on Digital Pathology* (pp. 83–90). Springer, Cham.
23. Weis C-A, Kather JN, Melchers S, Al-ahmdi H, Pollheimer MJ, Langner C, et al. Automatic evaluation of tumor budding in immunohistochemically stained colorectal carcinomas and correlation to clinical outcome. *Diagn Pathol*. 2018;13:64.